



Calhoun: The NPS Institutional Archive

Theses and Dissertations

Thesis Collection

1988-09

Small sample properties of bootstrap

Bernhardt, Stefan

<http://hdl.handle.net/10945/23394>



Calhoun is a project of the Dudley Knox Library at NPS, furthering the precepts and goals of open government and government transparency. All information contained herein has been approved for release by the NPS Public Affairs Officer.

Dudley Knox Library / Naval Postgraduate School
411 Dyer Road / 1 University Circle
Monterey, California USA 93943

<http://www.nps.edu/library>

NAVAL POSTGRADUATE SCHOOL

Monterey , California



THESIS

B452851

SMALL SAMPLE PROPERTIES OF BOOTSTRAP

by

Stefan Bernhardt

September 1988

Thesis Advisor

Toke Jayachandran

Approved for public release; distribution is unlimited

T238700

classified
by classification of this page

REPORT DOCUMENTATION PAGE				
Report Security Classification Unclassified		1b Restrictive Markings		
Security Classification Authority		3 Distribution Availability of Report		
Declassification Downgrading Schedule		Approved for public release; distribution is unlimited.		
Performing Organization Report Number(s)		5 Monitoring Organization Report Number(s)		
Name of Performing Organization Naval Postgraduate School		6b Office Symbol (if applicable) 55	7a Name of Monitoring Organization Naval Postgraduate School	
Address (city, state, and ZIP code) Monterey, CA 93943-5000		7b Address (city, state, and ZIP code) Monterey, CA 93943-5000		
Name of Funding Sponsoring Organization		8b Office Symbol (if applicable)	9 Procurement Instrument Identification Number	
Address (city, state, and ZIP code)		10 Source of Funding Numbers		
		Program Element No	Project No	Task No
				Work Unit Accession No
Title (include security classification) SMALL SAMPLE PROPERTIES OF BOOTSTRAP				
Personal Author(s) Stefan Bernhardt				
Type of Report Master's Thesis	13b Time Covered From To		14 Date of Report (year, month, day) September 1988	15 Page Count 63
Supplementary Notation The views expressed in this thesis are those of the author and do not reflect the official policy or position of the Department of Defense or the U.S. Government.				
Subject Codes		18 Subject Terms (continue on reverse if necessary and identify by block number)		
Group	Subgroup	Bootstrap, sampling distribution, confidence intervals		
Abstract (continue on reverse if necessary and identify by block number) The Bootstrap method is a nonparametric statistical technique for estimating the sampling distribution of estimators of unknown parameters. While the asymptotic theory for bootstrap is well established, this thesis investigates the behavior of the bootstrap for small sample sizes. For the exponential distribution and for normal linear regression the bootstrap estimates of the parameters and their variances are compared with the theoretical sampling distributions. The small sample properties of bootstrap confidence intervals using the percentile method and bias-corrected percentile method are also investigated.				
Distribution Availability of Abstract Unclassified unlimited <input type="checkbox"/> same as report <input type="checkbox"/> DTIC users		21 Abstract Security Classification Unclassified		
Name of Responsible Individual Rakesh Jayachandran		22b Telephone (include Area code) (408) 646-2600	22c Office Symbol 53Jy	

FORM 1473,84 MAR

83 APR edition may be used until exhausted
All other editions are obsolete

security classification of this page

Unclassified

Approved for public release; distribution is unlimited.

Small Sample Properties of Bootstrap

by

Stefan Bernhardt

Captain, Federal German Army

Dipl. Ing., Federal Army College Darmstadt, W. Germany, 1981

Submitted in partial fulfillment of the
requirements for the degree of

MASTER OF SCIENCE IN OPERATIONS RESEARCH

from the

NAVAL POSTGRADUATE SCHOOL

September 1988

ABSTRACT

The Bootstrap method is a nonparametric statistical technique for estimating the sampling distribution of estimators of unknown parameters. While the asymptotic theory for bootstrap is well established, this thesis investigates the behavior of the bootstrap for small sample sizes. For the exponential distribution and for normal linear regression the bootstrap estimates of the parameters and their variances are compared with the theoretical sampling distributions. The small sample properties of bootstrap confidence intervals using the percentile method and the bias-corrected percentile method are also investigated.

THESIS DISCLAIMER

The reader is cautioned that computer programs developed in this research may not have been exercised for all cases of interest. While every effort has been made, within the time available, to ensure that the programs are free of computational and logic errors, they cannot be considered validated. Any application of these programs without additional verification is at the risk of the user.

TABLE OF CONTENTS

I. INTRODUCTION	1
II. BOOTSTRAP METHODS - AN OVERVIEW	3
A. THE BASIC BOOTSTRAP METHOD	3
B. VARIATIONS OF THE BOOTSTRAP METHOD	5
1. Parametric Variations of Bootstrap	5
2. The Balanced Bootstrap	5
C. CONFIDENCE INTERVALS	6
1. The Percentile Method	6
2. The Bias-Corrected Percentile Method	6
III. THE EXPONENTIAL DISTRIBUTION	8
A. THEORETICAL RESULTS	8
B. THE SIMULATIONS	10
1. Point Estimation	10
a. Bias of the Bootstrap Estimate	11
b. Bootstrap Variance Estimation	13
2. Confidence Intervals	15
a. Simulation Validation	15
b. Coverage	16
c. Percentiles	18
IV. NORMAL LINEAR REGRESSION	20
A. THEORETICAL OVERVIEW	20
1. The Regression Model	20
2. Bootstrap Method for Regression Models	21
B. THE SIMULATIONS	22
1. Estimation of the Regression Coefficients	22
2. Bootstrap Estimates of the Variances	25
3. Confidence Intervals	25
4. Linear Regression with Mixtures of Normals	27

V. CONCLUSIONS	31
APPENDIX A. PERCENTILE ESTIMATION--PERCENTILE METHOD	33
APPENDIX B. PERCENTILE ESTIMATION--BIAS-CORRECTED PERCENTILE METHOD	35
APPENDIX C. VARIABILITY OF BOOTSTRAP POINT ESTIMATES--LINEAR REGRESSION	37
APPENDIX D. FORTRAN PROGRAM FOR BOOTSTRAP	39
APPENDIX E. SIMTBED DRIVER FOR BOOTSTRAP	46
LIST OF REFERENCES	52
INITIAL DISTRIBUTION LIST	54

LIST OF FIGURES

Figure 1. Probability Density Function of the Sampling Distribution	9
Figure 2. Average Bias	12
Figure 3. Bias of Bootstrap Variance Estimate	14

1. INTRODUCTION

The Bootstrap method, a statistical technique for estimating the sampling distributions of estimators of unknown parameters, was introduced by Efron [Ref. 1] in the mid 1970s. This computer intensive method is nonparametric in nature and relies on repeated resampling (bootstrapping) from the observed values of a random sample.

Suppose $x_1, x_2, x_3, \dots, x_n$ are the observed values of a random sample of size n , $X_1, X_2, X_3, \dots, X_n$, from a distribution $f_X(x; \theta)$. Let $\hat{\theta} = h(X_1, X_2, X_3, \dots, X_n)$ be an estimator for the unknown parameter θ . The sampling distribution of $\hat{\theta}$ completely describes the properties of the estimator and its knowledge would be useful for investigative purposes. However in many situations the analytical derivation of this distribution may be quite demanding. An alternative approach is to estimate the sampling distribution using bootstrap methods. A set of N bootstrap samples of size n , $x_{j1}^*, x_{j2}^*, x_{j3}^*, \dots, x_{jn}^*$ for $j = 1, 2, 3, \dots, N$ is generated by repeated uniform sampling with replacement from the set $\{x_1, x_2, x_3, \dots, x_n\}$. The estimate $\hat{\theta}_j^* = h(x_{j1}^*, x_{j2}^*, x_{j3}^*, \dots, x_{jn}^*)$ is computed for each of the N bootstrap samples. The empirical distribution of the $\hat{\theta}_j^*$ for $j = 1, 2, 3, \dots, N$ is taken as the estimate of the sampling distribution of $\hat{\theta}$.

Efron [Ref. 1] showed, that the bootstrap estimator is consistent and Beran et al. [Refs. 2, 3] proved that under fairly general regularity conditions the bootstrap distribution converges to the true sampling distribution as $n \rightarrow \infty$ and $N \rightarrow \infty$. It has also been demonstrated that bootstrap methods perform better than some of the other resampling techniques such as Hartigan's subsample method [Ref. 4] and the Tukey-Quenouille Jackknife [Ref. 1].

Although the asymptotic behavior of the bootstrap has been well established by theoretical research, there are still some problems dealing with the small sample properties of the methods, which are open for further investigation. One of these problems is the question of how the original sample size n and the number of bootstrap replications N affect the "closeness" of the bootstrap distribution to the exact sampling distribution. Another one deals with the applicability of bootstrap-based percentiles as a basis for estimating confidence intervals for parameters. Information about these issues will be useful to a practitioner in the decision of how to employ his resources.

The aim of this thesis is to address the two problems stated above. The approach which is taken is to consider probability distributions and their parameters, for which the

exact sampling distributions of the estimators can be derived theoretically. The results of simulations of the bootstrap method will be compared with the theoretical results in order to analyze the impact of the sample size n and the number of bootstrap replications N in the context of relatively small samples.

Chapter II provides an overview of some bootstrap methods and their properties. In Chapter III the bootstrap method is applied to the maximum likelihood estimator of the scale parameter of the exponential distribution. The estimation of the parameters in normal linear regression is studied in Chapter IV. In Chapter V the conclusions are presented.

II. BOOTSTRAP METHODS - AN OVERVIEW

A. THE BASIC BOOTSTRAP METHOD

The bootstrap method is a resampling technique for estimating the sampling distribution of an unknown parameter of a probability distribution. Let $\mathbf{X} = \{X_1, X_2, X_3, \dots, X_n\}$ be a sample of size n from a distribution with probability density function $f_x(x;\theta)$ and distribution function $F_x(x;\theta)$. Let $\hat{\theta} = h(\mathbf{X})$ be an estimator for the parameter θ . The distribution of $\hat{\theta}$, $g(\hat{\theta}; \theta)$ is called the sampling distribution of $\hat{\theta}$. In many problems it may be quite difficult to derive the sampling distribution analytically. But since computer resources are nowadays inexpensive and easily available, methods like bootstrap [Ref. 1], which will be described below, can be used to estimate the distribution of $\hat{\theta}$.

Suppose $\mathbf{x} = \{x_1, x_2, x_3, \dots, x_n\}$ are the observed values of the random sample. A bootstrap sample $\mathbf{x}^* = \{x^*_1, x^*_2, x^*_3, \dots, x^*_n\}$ ('*' indicates bootstrapping) of size n is obtained by randomly drawing with replacement from the original sample \mathbf{x} . Another way of describing this resampling procedure is: The empirical distribution function \hat{F} , which is discrete, is constructed by assigning a probability mass $1/n$ to each of the original samples x_i and then drawing n random samples from \hat{F} . Although it is possible to imagine, as Bickel and Freedman [Ref. 3] mention, bootstrap samples of an arbitrary size m , mathematical theory [Ref. 3] indicates that the use of the same size n as in the original sample is preferable.

Before continuing the description of the bootstrap method it seems appropriate to summarize some properties of any bootstrap sample. Each element in a bootstrap sample is drawn independently from the original sample. So conditional on the original sample the probability that the j th element in a bootstrap sample is any one of the original sample values is the same:

$$P\{X^*_j = x_i | \mathbf{x}\} = \frac{1}{n} \quad \text{for } i, j = 1, 2, 3, \dots, n. \quad (2.1)$$

The expectation, conditional on \mathbf{X} , of X^*_j is

$$E[X^*_j | \mathbf{x}] = \bar{x} \quad \text{for } j = 1, 2, 3, \dots, n, \quad (2.2)$$

where \bar{x} is the sample mean $\sum x_i/n$. Then for example the mean of the bootstrap sample $\sum x_i^*/n$ has the conditional expectation

$$E[\bar{X}^* | \mathbf{X}] = \bar{X} \quad (2.3)$$

and the unconditional expectation

$$E[\bar{X}^*] = E[E[\bar{X}^* | \mathbf{X}]] = \mu. \quad (2.4)$$

The variance of the mean of the bootstrap sample is

$$Var[\bar{X}^* | \mathbf{X}] = \frac{1}{n^2} \sum_{i=1}^n (x_i - \bar{x})^2 \quad (2.5)$$

which for $n \rightarrow \infty$ converges to $Var[\bar{X}]$.

The process of obtaining one bootstrap sample set and computing the estimator for this sample is called a bootstrap replication. For the bootstrap method N bootstrap replications are performed, where N varies throughout the literature between 100 and 2000. This means that N bootstrap samples \mathbf{x}_j^* for $j = 1, 2, 3, \dots, N$ are obtained and for each sample the estimator $\hat{\theta}_j^* = h(\mathbf{x}_j^*)$ is computed. The bootstrap distribution, the empirical distribution of the $\hat{\theta}_j^*$, is then an estimate of the sampling distribution of $\hat{\theta}$. The bootstrap estimate for θ is defined by

$$\hat{\theta}^* = \frac{1}{N} \sum_{j=1}^N \hat{\theta}_j^* \quad (2.6)$$

and

$$\hat{S}^* = \sqrt{\frac{1}{N-1} \sum_{j=1}^N (\hat{\theta}_j^* - \hat{\theta}^*)^2} \quad (2.7)$$

is the bootstrap estimate of $\sigma_{\hat{\theta}}$, the standard deviation of $\hat{\theta}$.

Efron [Ref. 1] and Bickel and Freedman [Ref. 3] have shown, that under fairly general regularity conditions, as $n \rightarrow \infty$ the bootstrap estimate and its standard deviation converge to their actual values.

B. VARIATIONS OF THE BOOTSTRAP METHOD

This section briefly describes some variations of the bootstrap method to demonstrate the variety of options available to the practitioner. These methods however will not be the subject of investigations in this thesis.

1. Parametric Variations of Bootstrap

To improve the bootstrap method in those cases, where additional information about the underlying distribution is available, Efron proposed [Ref. 4] the Smoothed Bootstrap. The major difference from the basic bootstrap is, that the bootstrap samples are now obtained by sampling from a continuous empirical distribution \tilde{F} . This distribution \tilde{F} is constructed by interpolating between the steps of the discrete empirical distribution \hat{F} using an appropriate smoothing function. Efron points out that the choice of the function is not arbitrary. In order to gain improvement of the results, compared to the basic bootstrap, the selection of the function type has to be compatible with the distribution under investigation. So this variation of the method is no longer nonparametric in an absolute sense.

If the exact distribution of the X_i is known except for the values of the parameters, this distribution can be used to perform the smoothing; Efron [Ref. 4] calls this method the Parametric Bootstrap.

2. The Balanced Bootstrap

Davison, Hinkley and Schechtman [Ref. 5] introduced the Balanced Bootstrap to eliminate the linear component of the bias of bootstrap estimators. Their method obtains the N bootstrap samples by first catenating the vector of n original samples N times, randomly permutating the resulting vector and then taking N successive vectors of size n , ensuring that each x_i occurs exactly N times in the total N bootstrap samples. It is easily seen, that when an estimator $h(\mathbf{X})$ for θ is linear and symmetric in \mathbf{X} , then

$$\frac{1}{N} \sum_{i=1}^N h(\mathbf{x}^*_i) = h(\mathbf{x}) . \quad (2.8)$$

C. CONFIDENCE INTERVALS

One of the applications of the sampling distribution is to approximate confidence intervals for a parameter. The following sections discuss two bootstrap-based methods for this purpose.

1. The Percentile Method

The percentile method is appealingly straightforward and provides, Efron [Ref. 4], good results. It is based on the definition of the empirical cumulative distribution function G^* of the estimator,

$$G^*(x) = \hat{P}\{\hat{\theta}^* \leq x\} = \frac{\#(\hat{\theta}_j^* \leq x)}{N} . \quad (2.9)$$

The p th percentile then can be approximated by $\hat{\theta}_p^*$, defined by

$$\hat{P}(\hat{\theta}^* \leq \hat{\theta}_p^*) \leq p . \quad (2.10)$$

Efron [Ref. 4] proposes the use of $(\hat{\theta}_\alpha^*, \hat{\theta}_{1-\alpha}^*)$ as an approximate $100(1 - 2\alpha)\%$ confidence interval for θ .

2. The Bias-Corrected Percentile Method

The bias-corrected percentile method covers those cases, where the empirical bootstrap distribution is not median-unbiased, i. e.,

$$P\{\hat{\theta}^* \leq \hat{\theta}\} \neq 0.5 . \quad (2.11)$$

The percentile method may produce inaccurate percentile estimates in this case. To compensate for these inaccuracies, Efron [Ref. 4] introduces the Bias-Corrected Percentile Method. This method relies, as Schenker [Ref. 6] points out, on an assumption, which in general is at best approximately valid. The assumption is, that there exists a function g such that

$$g(\hat{\theta}) - g(\theta) \sim N(\eta, \tau^2) \quad (2.12)$$

and

$$g(\hat{\theta}^*) - g(\hat{\theta}) \sim N(\eta, \tau^2) \quad (2.13)$$

with η and τ being real variables but constant for a specific case. Let

$$z_0 = \Phi^{-1}[G^*(\hat{\theta})] \quad (2.14)$$

and

$$z_x = \Phi^{-1}(1 - \alpha) \quad (2.15)$$

where Φ denotes the cumulative distribution function of the standard normal distribution and $\hat{\theta}$ is the value of the estimator for the original sample. Then the approximate $1 - 2\alpha$ confidence interval is given by

$$(G^{*-1}[\Phi(2z_0 - z_x)], G^{*-1}[\Phi(2z_0 + z_x)]) \quad (2.16)$$

It is easily seen, that for median unbiased sampling distributions, i. e., if

$$P\{\hat{\theta}^* \leq \hat{\theta}\} = 0.5, \quad (2.17)$$

$z_0 = 0$ and the bias-corrected percentile method is identical with the percentile method. Schenker's intention [Ref. 6] is to demonstrate some deficiencies of bootstrap-based confidence intervals for small sample sizes. Nevertheless, he does provide results which seem to indicate, that the bias-corrected percentile method is an improvement over the percentile method.

For the cases, where the underlying assumptions for the bias-corrected percentile method do not hold, Efron and Tibshirani [Ref. 7] proposed another method called the BC_x method. This thesis is concerned with the first two methods only.

III. THE EXPONENTIAL DISTRIBUTION

In this chapter, the performance of the Bootstrap method is compared to the theoretical results in the case where the underlying distribution is the exponential distribution.

A. THEORETICAL RESULTS

Let $X_1, X_2, X_3, \dots, X_n$ be i. i. d. random samples from the exponential distribution $\text{Exp}[\lambda]$ with probability density function

$$f_X(x) = \begin{cases} \lambda e^{-\lambda x} & \text{for } x > 0 \\ 0 & \text{otherwise} \end{cases}. \quad (3.18)$$

The Maximum Likelihood Estimator (MLE) for the scale parameter λ is

$$\hat{\lambda} = \frac{n}{\sum_{i=1}^n X_i}. \quad (3.19)$$

Using the fact, that the sum of n i.i.d. exponential random variables is distributed as Gamma $[\lambda, n]$, the probability density function of the random variable W , defined by

$$W = \frac{n}{\sum_{i=1}^n X_i}, \quad (3.20)$$

can be shown to be

$$f_W(w) = \begin{cases} \frac{(n\lambda)^n}{\Gamma(n)} \left(\frac{1}{w}\right)^{n+1} e^{-\frac{n\lambda}{w}} & \text{for } w > 0 \\ 0 & \text{otherwise} \end{cases}. \quad (3.21)$$

This is the exact sampling distribution of the maximum likelihood estimator for λ . Figure 1 shows the graph of the sampling distribution for $\lambda = 1$ and sample sizes $n = 10, 30$ and 50 .

Computations of the moments yield

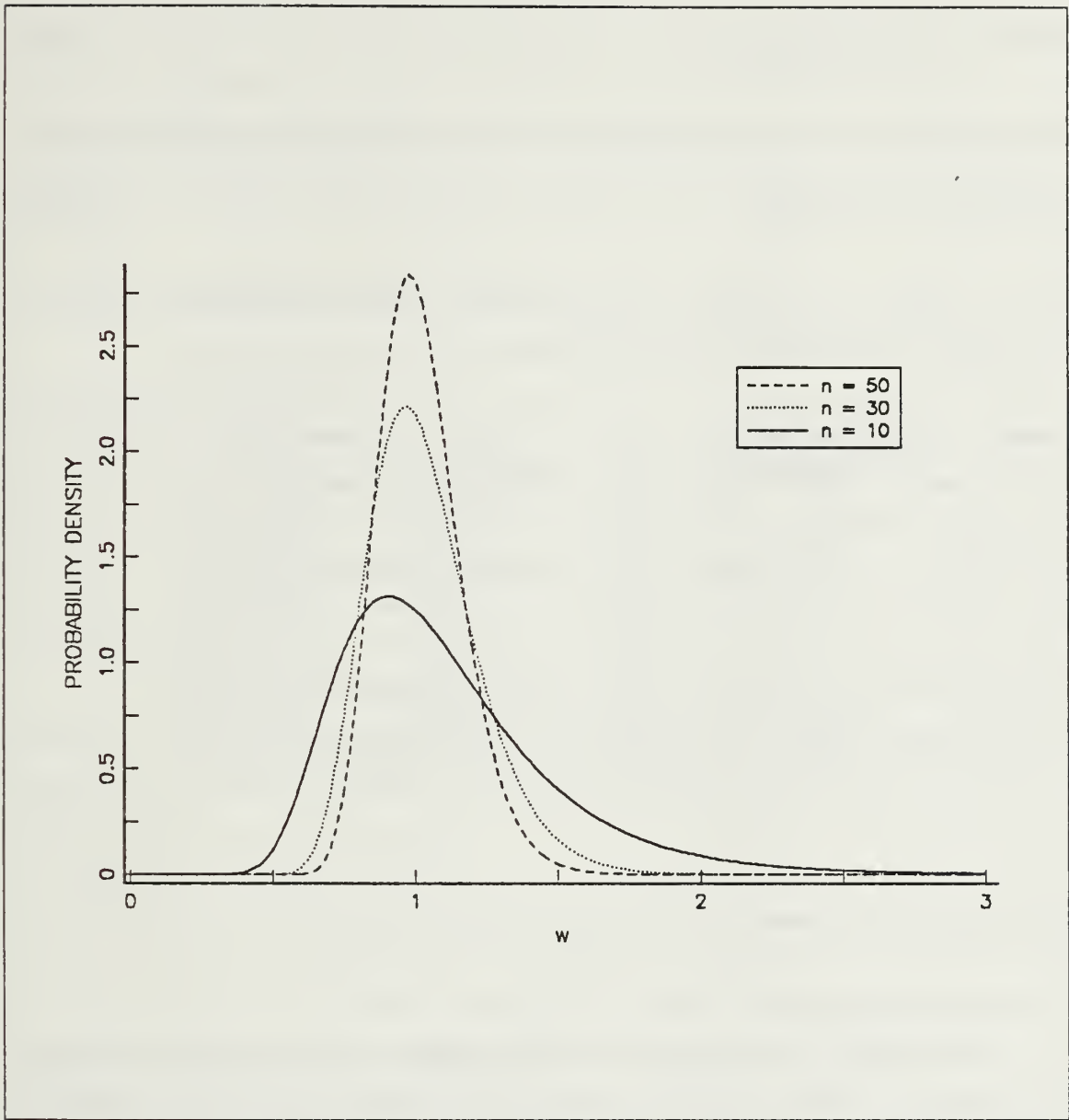


Figure 1. Probability Density Function of the Sampling Distribution: $\lambda = 1$, sample size $n = 10, 30, 50$.

$$E[W] = \frac{n}{n-1} \lambda, \quad (3.22)$$

which shows that the MLE is asymptotically unbiased, and

$$Var[W] = \frac{n^2}{(n-1)^2(n-2)} \lambda^2. \quad (3.23)$$

For this distribution exact probabilities can be computed using the following identity,

$$P\{W \leq w\} = 1 - I_n\left(\frac{n\lambda}{w}\right), \quad (3.24)$$

where I_n denotes the Incomplete Gamma function. Table 1 shows the true values for some percentiles of the distribution of W for $\lambda = 1$.

Table 1. PERCENTILES OF THE SAMPLING DISTRIBUTION: $\lambda = 1$, sample size n .

n	5%	10%	90%	95%
10	0.6367	0.7039	1.6074	1.8432
20	0.7174	0.7721	1.3769	1.5089
30	0.7587	0.8065	1.2915	1.3893
40	0.7852	0.8283	1.2446	1.3247
50	0.8042	0.8439	1.2142	1.2832
60	0.8187	0.8439	1.1926	1.2539

B. THE SIMULATIONS

1. Point Estimation

The purpose of this simulation is to investigate the performance of bootstrap point estimates. Cortes-Colon [Ref. 8] explored this subject for the sample mean of exponential variates, using the mean squared error as the criterion for his evaluation. This paper in contrast approaches the problem by looking at the bias and the variance separately in order to isolate effects.

The simulations in this section were conducted in SIMTBED [Ref. 9], a simulation software package for the IBM Personal Computer and compatibles, which uses a multiplicative congruential generator with multiplier 16807 and modulus $2^{31} - 1$ for the uniform and an acceptance-rejection scheme for the gamma variates. For the experiments in this section, ten super-replications were performed with differing numbers of trials for each original sample size. The original sample sizes used were $n = 10, 20, 30,$

40 and 50 with respective numbers of replications for one super-replication $M = 480, 240, 180, 120$ and 96 . With 10 super-replications, this sums up to a total of 4800, 2400, 1800, 1200 and 960 trials for each n and for each of the bootstrap replications. For validation purposes, similar simulations were performed on the author's personal computer using the APL language and also on the NPS mainframe using independent FORTRAN 77 programs. The results were similar to those obtained from SIMTBED.

a. Bias of the Bootstrap Estimate

In the first part of the simulation experiment, the quantity of interest is the bias B , the difference between the bootstrap estimate for the scale parameter λ and its true value ($\lambda = 1$). The bootstrap estimate $\hat{\lambda}^*$ was obtained according to equation 2.6 and the bias B was computed as $B = \hat{\lambda}^* - 1$ for each combination of n and N . Figure 2, created with GRAFSTAT [Ref. 10], shows the average values for B as a function of the number of bootstrap replications N , for various values of n . Table 2 lists the lengths of the central 90% confidence intervals for B , which are based on the super-replications.

The graph of the average values of B shows, that the number of bootstrap replications N has on the average almost no effect on the "closeness" of the bootstrap estimate to the actual value. Linear regression performed on the averages versus N resulted in slope parameters of the order of 10^{-3} and less. The bias is significantly affected by the sample size n . The reason for this behavior is the fact that the estimator is biased and that the bias decreases with increasing sample size. The average bias for each value of n is significantly larger than the amount expected from equation 3.22, which for this case would be $1/(n-1)$. The observed average bias is approximately twice the expected value which seems to indicate that the bootstrap method introduces additional bias. The variability of the bias as measured by the length of a 90% confidence interval is presented in Table 2. These lengths decrease with increasing sample size n but are not affected by the number of bootstrap replications N .

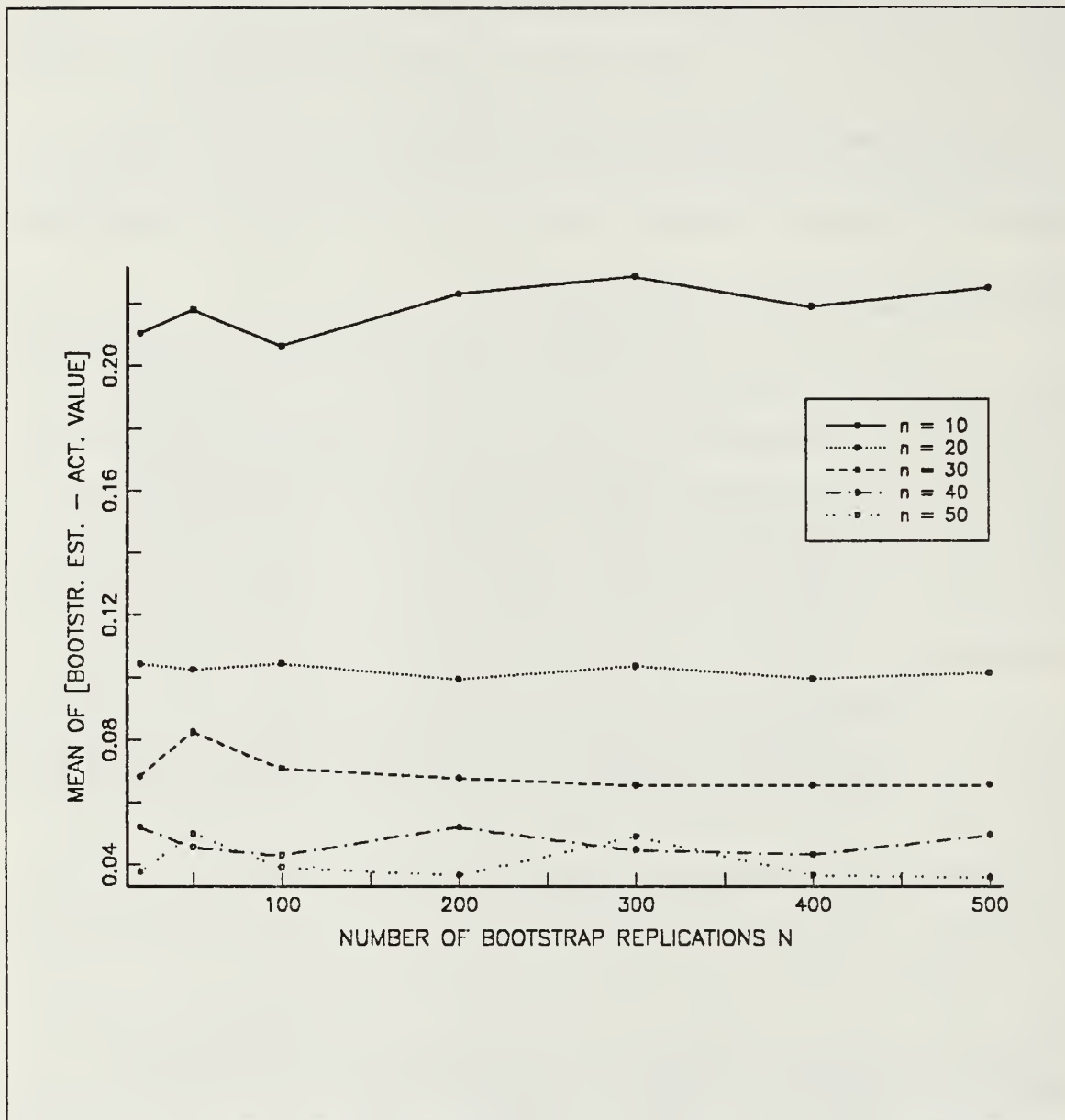


Figure 2. Average Bias: Average values of bootstrap estimate minus actual value for the true values $\lambda = 1$.

Table 2. VARIABILITY OF BIAS: Length of the 90% confidence interval for the bias B.
(Variances are less than 10^{-3})

N	n	10	20	30	40	50
20		1.3625	0.8567	0.6764	0.5804	0.5049
50		1.3718	0.8581	0.7064	0.5623	0.5100
100		1.3137	0.8275	0.6829	0.5518	0.4991
200		1.3481	0.8108	0.6694	0.5754	0.4665
300		1.3825	0.8476	0.6588	0.5692	0.5218
400		1.3574	0.8384	0.6839	0.5449	0.4996
500		1.3583	0.8512	0.6855	0.5742	0.5204

b. Bootstrap Variance Estimation

The quantity of interest here is the bias of the bootstrap estimate of the variance of $\hat{\lambda}$, i. e. $\hat{\sigma}^{*2} - \sigma^2$. The bootstrap estimate of the variance, $\hat{\sigma}^{*2}$, is computed according to equation 2.7 and σ^2 is the theoretical value from equation 3.23. The average values of the bias of the bootstrap variance estimate are displayed graphically in Figure 3 while the lengths of its 90% confidence intervals, depicting the variability, are listed in Table 3.

The graph shows that on the average bootstrap overestimates the variance of the maximum likelihood estimator of the scale parameter of the exponential distribution. The average bias after some fluctuation for low values of the bootstrap replications N seems to stabilize and from then on the number of bootstrap replications does not have a significant effect. Again the major impact on the bias is given by the sample size n . The graph clearly shows the decrease in bias with increasing n . The variability of the bias of the bootstrap variance estimate, represented by the lengths of the 90% confidence interval of the bias also does not seem to change with the number of bootstrap replications N . Least squares regression of the lengths on the number of bootstrap replications yields slope parameters of the order of 10^{-5} , which does not indicate a strong dependence. So a choice of about $N = 200$ bootstrap replications should be appropriate.

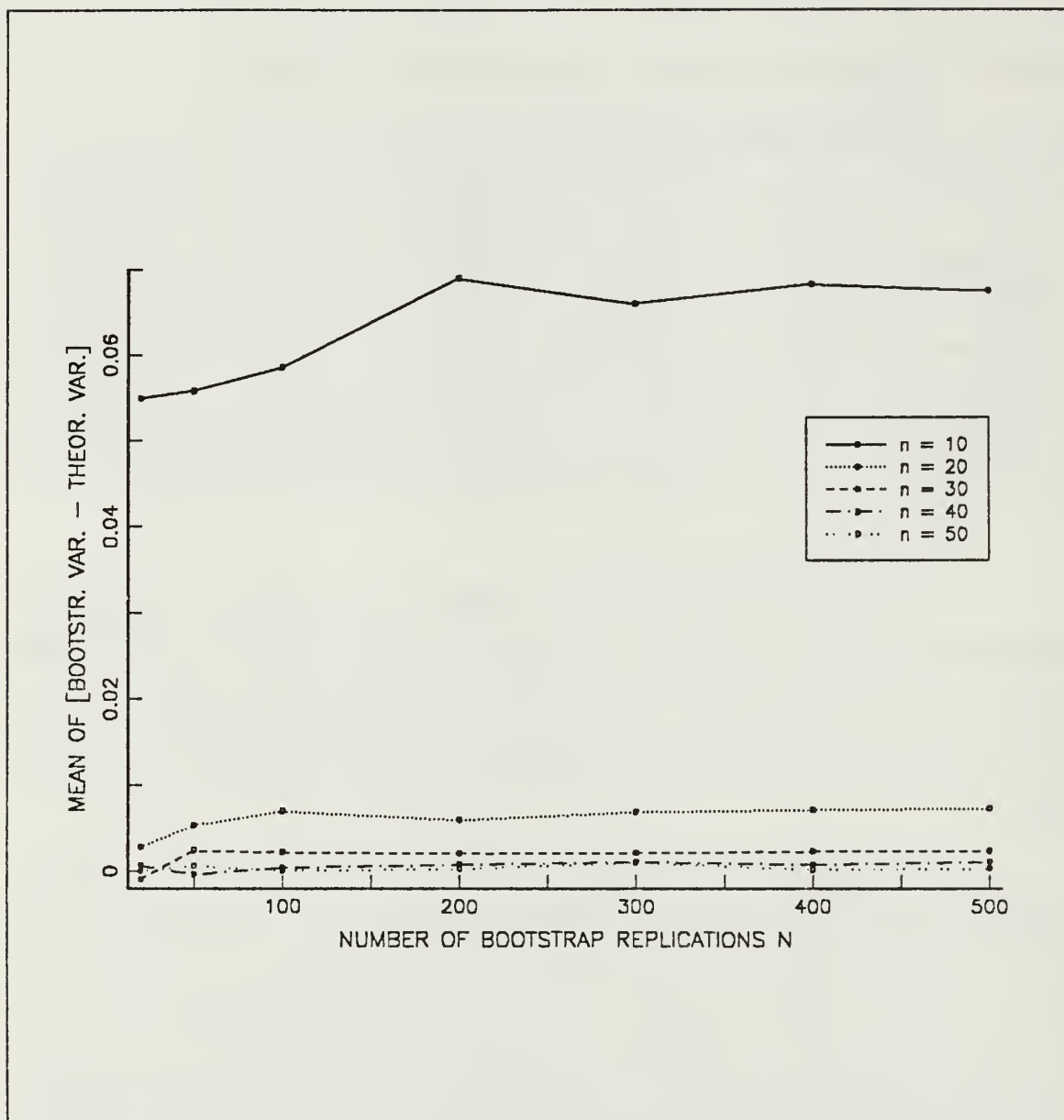


Figure 3. Bias of Bootstrap Variance Estimate: Average values of the bias of the bootstrap variance estimate $\hat{\sigma}^{*2} - \sigma^2$ for $\lambda = 1$.

Table 3. BIAS OF BOOTSTRAP VARIANCE ESTIMATE: Length of the 90% confidence interval of the bias of the bootstrap variance estimate $\hat{\sigma}^{*2} - \sigma^2$ for $\lambda = 1$.
(Variances are less than 10^{-3})

N	n	10	20	30	40	50
20		0.6988	0.1535	0.0742	0.0549	0.0383
50		0.6757	0.1568	0.0793	0.0456	0.0334
100		0.6266	0.1491	0.0725	0.0470	0.0310
200		0.7026	0.1452	0.0706	0.0436	0.0324
300		0.6806	0.1455	0.0718	0.0422	0.0328
400		0.6742	0.1469	0.0730	0.0450	0.0302
500		0.6572	0.1542	0.0710	0.0418	0.0297

The results of this section, briefly summarized, are: The number of bootstrap replications has no major impact on the "closeness" of the bootstrap estimates to the theoretical values. This observation is in agreement with the results by Cortes-Colon [Ref. 8] and by Efron and Tibshirani [Ref. 7].

2. Confidence Intervals

This section investigates bootstrap confidence intervals obtained by the percentile and bias-corrected percentile methods.

a. Simulation Validation

Validation is an important part of every simulation. Checking the results for plausibility, comparing them with the theory and with results obtained by other authors are some of the ways to accomplish validation. The latter way was specially chosen for this part of the thesis. To ensure that the percentile method and the bias-corrected percentile method were properly understood and correctly implemented in computer code, Efron's simulation [Ref. 4, page 84] was repeated. In the experiment random samples of size $n = 15$ are drawn from the exponential distribution $\text{Exp}[\lambda = 1]$. The sample are standardized to ensure that the sample mean $\bar{x} = 0$ and the sample variance $\sum(x_i - \bar{x})^2/(n - 1) = 1$. The bootstrap method is then applied to the standardized samples with the number of bootstrap replications $N = 1000$. Selected percentiles are approximated using the percentile method and the bias-corrected percentile method. Table

4 shows the results for 10 trials. The averages of the estimated percentiles over the ten trials and the corresponding results obtained by Efron [Ref. 4, p. 85] are also presented. The numbers obtained are quite close to Efron's results. The simulation was programmed in FORTRAN 77 and conducted on the NPS IBM mainframe. The random variates, exponential and uniform, were generated using the random number package LLRANDOMII [Ref. 11]. Appendix D shows the listing of the program for the percentile method and the bias-corrected percentile method.

Table 4. SIMULATION VALIDATION: Nonparametric confidence intervals of exponential variates $\text{Exp}[\lambda = 1]$, standardized, i. e. sample mean = 0 and sample variance = 1; $n = 15$, $N = 1000$.

Trial	Percentile Method				Bias-corrected PM			
	5%	10%	90%	95%	5%	10%	90%	95%
1	-0.358	-0.300	0.368	0.457	-0.358	-0.300	0.368	0.457
2	-0.403	-0.322	0.324	0.425	-0.384	-0.293	0.359	0.454
3	-0.377	-0.298	0.305	0.435	-0.373	-0.290	0.328	0.453
4	-0.375	-0.309	0.331	0.433	-0.373	-0.304	0.340	0.440
5	-0.381	-0.300	0.329	0.431	-0.378	-0.292	0.343	0.439
6	-0.408	-0.302	0.345	0.451	-0.391	-0.288	0.355	0.463
7	-0.347	-0.302	0.330	0.478	-0.322	-0.271	0.399	0.565
8	-0.391	-0.304	0.320	0.426	-0.356	-0.289	0.348	0.449
9	-0.384	-0.320	0.309	0.410	-0.371	-0.298	0.336	0.442
10	-0.425	-0.332	0.332	0.404	-0.401	-0.305	0.362	0.436
Average	-0.385	-0.309	0.329	0.435	-0.371	-0.293	0.354	0.460
Efron	-0.39	-0.32	0.33	0.43	-0.36	-0.29	0.36	0.47

b. Coverage

The interpretation of a confidence interval, e. g. 90%, for a parameter of interest is, that in the long run with a relative frequency of 0.9, the computed confidence intervals cover the actual value of the parameter. Thus the relative frequency of coverage can be used to assess the quality and applicability of a method, which produces confidence intervals. In this section, the coverage is investigated for the percentile method and the bias-corrected percentile method.

The simulation looks at the central 90% confidence interval. This interval is set up using the 5th and 95th quantiles of the empirical bootstrap distribution for the scale parameter λ of the exponential distribution for both methods. The simulation was programmed in FORTRAN 77 and run on the NPS mainframe computer. Random numbers were generated with LLRANDOMII [Ref. 11]. For each combination of

sample size n and number of bootstrap replications N the simulation consists of 1000 repetitions, for each of which the coverage of the actual value $\lambda = 1$ was checked. Table 5 shows the counts for the percentile method and Table 6 for the bias-corrected percentile method.

Table 5. COVERAGE--PERCENTILE METHOD CONFIDENCE INTERVAL: Coverage of the true value $\lambda = 1$ by the 90% confidence interval obtained from the percentile method, out of 1000 repetitions.

N	n	10	20	30	40	50
50		826	845	859	883	888
100		804	851	877	888	889
200		794	845	881	853	881
300		819	856	862	870	884
500		784	848	840	876	860

Table 6. COVERAGE--BIAS-CORRECTED PERCENTILE METHOD CONFIDENCE INTERVAL: Coverage of the true value $\lambda = 1$ by the 90% confidence interval obtained from the bias-corrected percentile method, out of 1000 repetitions.

N	n	10	20	30	40	50
50		831	831	855	875	870
100		815	847	880	888	888
200		792	851	881	859	873
300		821	871	857	871	888
500		793	850	847	883	860

The coverage in all cases is below the nominal level of 90%. The coverage appears to be somewhat erratic for the smaller values of sample sizes, $n = 10$ and 20, but it seems to improve with increasing n . Schenker [Ref. 6] observed a similar behavior in his investigation dealing with the estimation of the variance of a normal distribution. The number of bootstrap replications N again seems not to have a significant effect. Significant differences between the percentile method and the bias-corrected percentile method are also not detectable in this experiment.

c. Percentiles

The simulation in the previous section was set up to also provide the average values of the 5th, 10th, 90th and 95th percentile of the empirical bootstrap distribution. Table 7 lists these values for $N = 500$ bootstrap replications; these are averages of 1000 trials. Both methods, percentile and bias-corrected percentile method on the average overestimate the percentiles compared to the theoretical values from Table 1. The amount of overestimation is shown in the table in parentheses. This amount is in general larger for the percentile method than for the bias-corrected percentile method, which means that the correction, which the latter method applies, is working in the right direction. The difference between theoretical values and the bootstrap-based estimates decreases with increasing original sample size n .

Table 7. AVERAGE PERCENTILES: Average values for percentiles obtained with the percentile and the bias-corrected percentile method in 1000 trials; $\lambda = 1$, number of bootstrap replications $N = 500$; numbers in parentheses are the amount of overestimation, compared to the theoretical values.

n	Percentile Method				Bias-corrected PM			
	5%	10%	90%	95%	5%	10%	90%	95%
10	0.755 (0.118)	0.817 (0.113)	1.727 (0.129)	1.979 (0.136)	0.737 (0.100)	0.795 (0.091)	1.652 (0.045)	1.890 (0.047)
20	0.773 (0.056)	0.825 (0.053)	1.426 (0.049)	1.560 (0.051)	0.757 (0.040)	0.808 (0.036)	1.389 (0.012)	1.517 (0.008)
30	0.799 (0.040)	0.846 (0.039)	1.327 (0.035)	1.424 (0.035)	0.785 (0.026)	0.831 (0.024)	1.300 (0.008)	1.395 (0.006)
40	0.814 (0.010)	0.857 (0.013)	1.270 (0.025)	1.350 (0.025)	0.802 (-0.002)	0.844 (0.000)	1.250 (0.005)	1.328 (0.003)
50	0.833 (0.029)	0.872 (0.028)	1.240 (0.026)	1.308 (0.025)	0.824 (0.020)	0.863 (0.019)	1.225 (0.011)	1.292 (0.009)

The behavior of percentile estimates was investigated further. The simulation for this purpose was done in SIMTBED [Ref. 9] on the author's personal computer. The 5th and 95th percentiles were selected as representative objects for investigation. The number of trials is 1200 for $n = 10$, 600 for $n = 20$, 480 for $n = 30$, 300 for $n = 40$ and 240 for $n = 50$. Appendix A lists the results for the percentile method. These results show that both the standard deviation of the percentile estimate and the width of the central 90% confidence interval decrease with increasing sample

size n . The number of bootstrap replications N seems not to affect the results. Least squares regression of the values on the number of bootstrap replications resulted in values for the slope of 10^{-4} and less. And tests for distributional fit in GRAFSTAT [Ref. 10] did not show significant differences for different numbers of bootstrap replications.

The simulation was repeated for the 5th and 95th percentiles using the bias-corrected percentile method. The results are listed in Appendix B. The conclusions for this method are basically the same as with the percentile method, decreasing standard deviation and width of the central 90% confidence interval with increasing sample size and no effect of the number of bootstrap replications. The only difference again is that the bias-corrected percentile method is on the average closer to the theoretical value than the percentile method.

IV. NORMAL LINEAR REGRESSION

The results of simulations to study the properties of bootstrap estimators of the parameters in a simple linear regression model are presented in this chapter.

A. THEORETICAL OVERVIEW

1. The Regression Model

Let $\{(x_1, y_1), (x_2, y_2), (x_3, y_3), \dots, (x_n, y_n)\}$ be n pairs of observations with x as the independent variable and y the dependent variable. Under the assumptions of independence, normal distribution and homoscedasticity for the random variables Y_i , the model for a linear relationship between x and Y_i is

$$Y_i = \beta_0 + \beta_1 x + \varepsilon_i \quad \text{for } i = 1, 2, 3, \dots, n. \quad (4.25)$$

The random variables ε_i have mean 0 and variance σ^2 and are normally distributed:

$$\varepsilon_i \sim N(0, \sigma^2) \quad \text{for } i = 1, 2, 3, \dots, n. \quad (4.26)$$

It is well known that the maximum likelihood estimates for the coefficients β_0 and β_1 are

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n x_i y_i - n \bar{x} \bar{y}}{\sum_{i=1}^n x_i^2 - n \bar{x}^2} \quad (4.27)$$

and

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}. \quad (4.28)$$

Both estimators are unbiased, i. e.

$$E[\hat{\beta}_0] = \beta_0 \quad \text{and} \quad E[\hat{\beta}_1] = \beta_1. \quad (4.29)$$

The joint sampling distribution for $\hat{\beta}_0$ and $\hat{\beta}_1$ is known to be a bivariate normal distribution. The marginal distributions are normal with means equal to the respective true values and the variances

$$Var[\hat{\beta}_0] = \frac{\sigma^2}{n} \frac{\sum_{i=1}^n x_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (4.30)$$

and

$$Var[\hat{\beta}_1] = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} . \quad (4.31)$$

The covariance between $\hat{\beta}_0$ and $\hat{\beta}_1$ is

$$Cov[\hat{\beta}_0, \hat{\beta}_1] = - \frac{\sigma^2 \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2} . \quad (4.32)$$

2. Bootstrap Method for Regression Models

The implementation of the bootstrap method for regression models [Ref. 4] differs slightly from the one in the one-sample case. It is described here for the normal linear regression of one dependent and one independent variable, which is the topic of this chapter.

To perform the bootstrap, first the least squares estimates $\hat{\beta}_0$ and $\hat{\beta}_1$ (equations 4.27 and 4.28) are computed. These estimates are used to compute the residuals e_i :

$$e_i = y_i - (\hat{\beta}_0 + \hat{\beta}_1 x_i) \text{ for } i = 1, 2, 3, \dots, n. \quad (4.33)$$

A bootstrap sample \mathbf{e}^* of size n , which is of the same size as the original sample, is obtained by randomly drawing with replacement n times from the e_i . Computing

$$y_i^* = \hat{\beta}_0 + \hat{\beta}_1 x_i + e_i^* \text{ for } i = 1, 2, 3, \dots, n \quad (4.34)$$

results in n pairs of 'observations' $\{(x_1, y_1^*), (x_2, y_2^*), (x_3, y_3^*), \dots, (x_n, y_n^*)\}$. These n pairs of 'observations' are used to compute the bootstrap estimates $\hat{\beta}_0^*$ and $\hat{\beta}_1^*$ using the equations 4.27 and 4.28. The process of randomly drawing and computing the estimates is repeated for a total of N bootstrap replications. The bootstrap estimates $\hat{\beta}_{0,j}^*$ and

$\hat{\beta}_{1,j}^*$ for $j = 1, 2, 3, \dots, N$ are used to construct the empirical sampling distribution for the estimators $\hat{\beta}_0$ and $\hat{\beta}_1$. All other quantities then can be estimated as described in Chapter II.

B. THE SIMULATIONS

The choice of the values for x , the independent variable in the regression model influences the variability of the results. Since this effect was not to be investigated, the values for x were kept fixed throughout the simulations. The values for the x were evenly spread from $10/n$ to 10 in increments of $10/n$ where n indicates the sample size. The values chosen for the coefficients β_0 and β_1 and the variance σ^2 are discussed below. For each simulation the sample pairs (x_i, y_i) were obtained by computing

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i \text{ for } i = 1, 2, 3, \dots, n, \quad (4.35)$$

with the ε_i being normal random variates with mean zero and variance σ^2 . Having set up the pairs of observations, the simulation then proceeds as described in the previous section.

1. Estimation of the Regression Coefficients

The first simulation estimated the coefficients β_0 and β_1 by applying equations 4.27 and 4.28 to each bootstrap sample and averaging over N , the number of bootstrap replications. It was conducted in SIMTBED [Ref. 9] on the author's personal computer and consisted of ten super replications. For $n = 10, 20, 30, 40$ and 50 the respective number of trials was $300, 150, 100, 75$ and 60 within each super replication. This sums up to a total number of trials of 3000 for $n = 10$, 1500 for $n = 20$, 1000 for $n = 30$ etc.. A simulation with those ten super replications was conducted for each selected number of bootstrap replications N . For this simulation the values of the parameters were $\beta_0 = 1.5$, $\beta_1 = 0.8$ and $\sigma^2 = 0.5$. The tables below show the averages of the estimates over all super replications, Table 8 for $\hat{\beta}_0^*$ and Table 9 for $\hat{\beta}_1^*$.

Table 8. **BOOTSTRAP ESTIMATE: Y-INTERCEPT:** Average values and standard deviations () of the bootstrap estimate of the y-intercept; actual value $\beta_0 = 1.5$.

N	n	10	20	30	40	50
50		1.507 (0.4782)	1.499 (0.3348)	1.487 (0.2657)	1.505 (0.2286)	1.492 (0.2073)
100		1.507 (0.4835)	1.501 (0.3253)	1.491 (0.2692)	1.500 (0.2264)	1.514 (0.2043)
200		1.497 (0.4777)	1.485 (0.3224)	1.492 (0.2588)	1.498 (0.2283)	1.506 (0.2017)
300		1.489 (0.4848)	1.506 (0.3390)	1.504 (0.2586)	1.493 (0.2316)	1.496 (0.1961)
400		1.492 (0.4799)	1.519 (0.3226)	1.482 (0.2715)	1.493 (0.2288)	1.509 (0.2054)
500		1.507 (0.4838)	1.499 (0.3251)	1.509 (0.2605)	1.508 (0.2292)	1.504 (0.1999)

Table 9. **BOOTSTRAP ESTIMATE: SLOPE:** Average values and standard deviations () of the bootstrap estimate of the slope parameter; actual value $\beta_1 = 0.8$.

N	n	10	20	30	40	50
50		0.7991 (0.07755)	0.8006 (0.05575)	0.8011 (0.04460)	0.7990 (0.03988)	0.8006 (0.03353)
100		0.7981 (0.07837)	0.7997 (0.05456)	0.8011 (0.04596)	0.8002 (0.03935)	0.7985 (0.03502)
200		0.8003 (0.07665)	0.8029 (0.05315)	0.8015 (0.04503)	0.8012 (0.03845)	0.7985 (0.03538)
300		0.8002 (0.07778)	0.7986 (0.05561)	0.8006 (0.04452)	0.8011 (0.03958)	0.7992 (0.03383)
400		0.8008 (0.07835)	0.7969 (0.05486)	0.8027 (0.04690)	0.8017 (0.03950)	0.7986 (0.03563)
500		0.7989 (0.07765)	0.8008 (0.05451)	0.7987 (0.04395)	0.7987 (0.03948)	0.7995 (0.03426)

It can be seen that on the average the bootstrap estimates of the regression parameters are fairly close to the theoretical values. The number of bootstrap replications

N again seems not to affect the "closeness" of the bootstrap estimates to the theoretical values. The computed standard deviations of the estimates (shown in parentheses) also confirm this conclusion. As a more detailed representation of the variability, Appendix C shows selected quantiles of the bootstrap estimates of the regression coefficients from the simulation with $N = 300$. The simulation results compare favorably with the theory.

For selected numbers of bootstrap replications the simulation was repeated with different sets of values for the parameters β_0 , β_1 and σ^2 . The following tables show the simulation results.

Table 10. ESTIMATION OF THE REGRESSION COEFFICIENTS: Average values and standard deviations of the bootstrap estimates of the regression coefficients; $N = 300$; theoretical values: $\beta_0 = 0.5$, $\beta_1 = 2.0$, $\sigma^2 = 0.5$.

n	10	20	30	40	50
$\hat{\beta}_0^*$	0.4893 (0.4848)	0.5055 (0.3390)	0.5041 (0.2586)	0.4929 (0.2316)	0.4957 (0.1961)
$\hat{\beta}_1^*$	2.000 (0.0778)	1.999 (0.0556)	2.001 (0.0445)	2.001 (0.0396)	1.999 (0.0338)

Table 11. ESTIMATION OF THE REGRESSION COEFFICIENTS: Average values and standard deviations of the bootstrap estimates of the regression coefficients; $N = 200$; theoretical values: $\beta_0 = 1.5$, $\beta_1 = 0.8$, $\sigma^2 = 4.0$.

n	10	20	30	40	50
$\hat{\beta}_0^*$	1.492 (1.351)	1.457 (0.9119)	1.477 (0.7319)	1.494 (0.6458)	1.516 (0.5705)
$\hat{\beta}_1^*$	0.8009 (0.2168)	0.8082 (0.1503)	0.8042 (0.1274)	0.8033 (0.1087)	0.7958 (0.1001)

The averages of the bootstrap estimates are again quite close to the actual values. The changes in the standard deviations clearly reflect the changes in the parameter values and conform with the theory.

2. Bootstrap Estimates of the Variances

The same setup, as far as the number of repetitions and actual values for the parameters are concerned, was used for this simulation. The quantities under investigation now were the differences between the bootstrap estimates and the theoretical values of the variances and the covariance of the least squares estimators. The bootstrap estimates for the variances and the covariance were obtained following equation 2.7, while the theoretical values were computed according to equations 4.30, 4.31 and 4.32. The average differences for the variances turned out to be negative, i. e. the bootstrap estimate is on the average lower than the theoretical value. For the covariance which is negative the average differences were positive which means that the absolute value of the bootstrap covariance estimate on the average is lower than the theoretical value. The average absolute values of the differences were less than 0.05 and decreasing with increasing original sample size n . The number of bootstrap replications N had no effect on the average difference. Standard deviations for the differences were of the order 0.1, decreasing with increasing n .

Again for some selected cases the simulation was repeated with different sets of values for the parameters β_0 , β_1 and σ^2 . The following Table 12 shows the results of one of these.

Table 12. BOOTSTRAP VARIANCE ESTIMATE: Average values and standard deviations of the difference $\Delta\text{Var}^*[\hat{\beta}]$ between the bootstrap variance estimates of the regression coefficients and the theoretical value; $N = 200$; theoretical values: $\beta_0 = 1.5$, $\beta_1 = 0.8$, $\sigma^2 = 4.0$.

n	10	20	30	40	50
$\Delta\text{Var}^*[\hat{\beta}_0]$	-0.3605 (0.7596)	-0.0921 (0.2765)	-0.0310 (0.1509)	-0.0271 (0.0970)	-0.0124 (0.0732)

The bootstrap estimates of the variance of $\hat{\beta}_1$ and the covariance between $\hat{\beta}_0$ and $\hat{\beta}_1$ were also quite accurate; the differences from the corresponding theoretical values were of the order 0.05 with standard deviations of the order 0.1.

3. Confidence Intervals

In this section the bootstrap estimates for percentiles as bounds for confidence intervals for the estimates of the regression coefficients are investigated. The investi-

gation focuses on the coverage of the central 90% confidence interval using both the percentile and the bias-corrected percentile method. The simulation was written in FORTRAN 77 for the IBM mainframe computer. It used the theoretical values $\beta_0 = 1.5$, $\beta_1 = 0.8$ and $\sigma^2 = 0.5$. 500 trials were used for each selected combination of sample size n and bootstrap replications N . Tables 13 and 14 lists the results for the percentile method and Tables 15 and 16 for the bias-corrected percentile method. The coverage for both methods is in most cases below the nominal 90% level although the differences are fairly small. The increase in coverage with increasing sample size n seems obvious while the number of bootstrap replications N does not seem to have any influence. A significant difference between the two methods can not be demonstrated with the results.

Table 13. COVERAGE - PERCENTILE METHOD CONFIDENCE INTERVAL: Percentage of coverage of the true value $\beta_0 = 1.5$ by the bootstrap 90% confidence interval obtained from the percentile method, in 500 trials.

N	n	10	20	30	40	50
100		81.0	82.2	88.0	88.4	88.6
200		82.2	85.6	87.2	88.8	87.2
300		83.2	86.2	84.2	83.4	90.0
400		83.0	84.0	90.6	89.8	88.4
500		82.8	85.4	88.6	87.8	88.0

Table 14. COVERAGE - PERCENTILE METHOD CONFIDENCE INTERVAL: Percentage of coverage of the true value $\beta_1 = 0.8$ by the bootstrap 90% confidence interval obtained from the percentile method, in 500 trials.

N	n	10	20	30	40	50
100		80.4	85.0	87.6	89.6	88.2
200		80.4	84.8	88.0	88.4	89.6
300		82.6	83.6	83.0	85.0	89.0
400		80.4	84.2	88.4	88.8	88.0
500		81.4	85.2	89.0	87.4	88.2

Table 15. **COVERAGE--BIAS-CORRECTED PERCENTILE METHOD CONFIDENCE INTERVAL:** Percentage of coverage of the true value $\beta_0 = 1.5$ by the bootstrap 90% confidence interval obtained from the bias-corrected percentile method, out of 500 repetitions.

N	n	10	20	30	40	50
100		80.6	83.2	87.6	87.8	87.8
200		82.0	86.8	86.4	88.0	86.6
300		83.2	85.6	84.4	83.8	90.4
400		83.6	84.2	89.0	89.6	87.8
500		82.8	84.2	89.6	88.0	88.2

Table 16. **COVERAGE--BIAS-CORRECTED PERCENTILE METHOD CONFIDENCE INTERVAL:** Percentage of coverage of the true value $\beta_1 = 0.8$ by the bootstrap 90% confidence interval obtained from the bias-corrected percentile method, in 500 trials.

N	n	10	20	30	40	50
100		79.8	86.0	89.0	88.2	88.4
200		80.2	85.4	88.4	88.6	89.2
300		83.6	84.2	83.8	85.0	88.8
400		79.2	84.4	89.0	88.2	87.8
500		81.4	86.0	88.2	86.2	87.6

4. Linear Regression with Mixtures of Normals

As a further test case for the bootstrap of regression models, another linear regression model was chosen. In this model all the assumptions about the dependent random variable Y as in the previous model hold, except the distributional assumption. Now the underlying distribution is assumed to be a mixture of two normal distributions. The model equation 4.25 still holds, but now equation 4.26 becomes

$$\varepsilon_i \sim \begin{cases} N(\mu_1, \sigma_1^2) & \text{with probability } p \\ N(\mu_2, \sigma_2^2) & \text{with probability } (1-p) \end{cases} \quad \text{for } i = 1, 2, 3, \dots, n. \quad (4.36)$$

The expectation for ε is

$$E[\varepsilon] = p\mu_1 + (1 - p)\mu_2. \quad (4.37)$$

In the simulations $E[\varepsilon]$ is for convenience set equal to zero by adjusting the values of p , μ_1 and μ_2 appropriately. The resulting variance of the mixture of two normal distributions is

$$\sigma_{res}^2 = p\sigma_1^2 + (1 - p)\sigma_2^2 + p(1 - p)(\mu_1 - \mu_2)^2. \quad (4.38)$$

The same simulation setup as for the normal linear regression in SIMTBED was used to conduct simulations for this regression model. The following tables, Table 17 and Table 18 show the results of two simulation runs.

Table 17. ESTIMATION OF THE REGRESSION COEFFICIENTS: Average values and standard deviations of the bootstrap estimates of the regression coefficients; $N = 100$; theoretical values: $\beta_0 = 1.5$, $\beta_1 = 0.8$, $p = 0.5$, $\mu_1 = 1.5$, $\sigma_1^2 = 0.5$, $\mu_2 = -1.5$, $\sigma_2^2 = 0.5$.

n	10	20	30	40	50
$\hat{\beta}_0^*$	1.488 (1.137)	1.511 (0.7922)	1.498 (0.6319)	1.490 (0.5180)	1.475 (0.4981)
$\hat{\beta}_1^*$	0.8052 (0.1824)	0.8014 (0.1317)	0.8003 (0.1081)	0.8034 (0.0885)	0.8048 (0.0834)

Table 18. ESTIMATION OF THE REGRESSION COEFFICIENTS: Average values and standard deviations of the bootstrap estimates of the regression coefficients; $N = 200$; theoretical values: $\beta_0 = 1.5$, $\beta_1 = 0.8$, $p = 0.25$, $\mu_1 = -1.5$, $\sigma_1^2 = 0.5$, $\mu_2 = 0.5$, $\sigma_2^2 = 1.0$.

n	10	20	30	40	50
$\hat{\beta}_0^*$	1.508 (0.8713)	1.512 (0.5979)	1.512 (0.4638)	1.490 (0.4145)	1.500 (0.3660)
$\hat{\beta}_1^*$	0.7990 (0.1397)	0.7968 (0.0973)	0.7980 (0.0798)	0.8053 (0.0712)	0.7971 (0.0620)

The estimates for the parameters are once more close to the actual values.

For $\beta_0 = 1.5$, $\beta_1 = 0.8$, $p = 0.25$, $\mu_1 = -1.5$, $\sigma_1^2 = 0.5$, $\mu_2 = 0.5$ and $\sigma_2^2 = 1.0$ 90% bootstrap confidence intervals were computed. The simulation consisted of 500 trials for each combination of the sample size n and the bootstrap replications N . The following four tables, Table 19 through Table 22 contain the results of the simulation,

Tables 19 and 20 for the percentile method and Tables 21 and 22 for the bias corrected percentile method.

Table 19. COVERAGE - PERCENTILE METHOD CONFIDENCE INTERVAL: Percentage of coverage of the true value $\beta_0 = 1.5$ by the bootstrap 90% confidence interval obtained from the percentile method, in 500 trials.

N	n	10	20	30	40	50
100		78.4	84.0	88.2	88.4	88.6
200		80.4	85.8	87.6	88.2	86.6
300		83.6	84.8	84.8	85.0	87.8
400		81.4	85.0	88.0	91.0	86.6
500		84.6	84.6	87.6	88.4	88.2

Table 20. COVERAGE - PERCENTILE METHOD CONFIDENCE INTERVAL: Percentage of coverage of the true value $\beta_1 = 0.8$ by the bootstrap 90% confidence interval obtained from the percentile method, in 500 trials.

N	n	10	20	30	40	50
100		80.2	86.4	88.4	88.2	87.2
200		81.0	83.6	89.6	87.4	88.4
300		82.6	87.2	84.8	86.6	89.4
400		82.6	87.8	86.6	89.8	87.8
500		81.0	85.4	87.8	87.8	88.2

The results are almost the same as for the standard regression model. The coverages are with two exceptions below the nominal level; they increase with increasing sample size n and the bootstrap replications N again seem not to affect the coverage. Also no significant difference is detectable between the coverages of the confidence intervals from the percentile method and the bias-corrected percentile method.

Table 21. COVERAGE--BIAS-CORRECTED PERCENTILE METHOD CONFIDENCE INTERVAL: Percentage of coverage of the true value $\beta_0 = 1.5$ by the bootstrap 90% confidence interval obtained from the bias-corrected percentile method, out of 500 trials.

N	n	10	20	30	40	50
100		80.0	83.6	88.2	87.6	87.8
200		81.6	84.4	86.6	87.2	86.2
300		83.0	85.0	85.6	85.4	89.0
400		83.4	85.0	87.6	91.0	86.4
500		82.0	84.8	88.0	88.6	87.6

Table 22. COVERAGE--BIAS-CORRECTED PERCENTILE METHOD CONFIDENCE INTERVAL: Percentage of coverage of the true value $\beta_1 = 0.8$ by the bootstrap 90% confidence interval obtained from the bias-corrected percentile method, in 500 trials.

N	n	10	20	30	40	50
100		80.8	84.8	87.6	88.0	87.4
200		81.2	83.8	88.6	87.4	86.6
300		82.2	86.4	84.8	86.0	89.2
400		82.8	87.4	86.4	89.2	87.4
500		80.2	84.4	88.4	87.8	88.2

V. CONCLUSIONS

The subject of this thesis is an investigation of the performance of the bootstrap method for small sample sizes in the two scenarios of the exponential distribution and linear regression. In both cases the simulation results show that the bootstrap method provides reasonable approximations for the estimation of statistical parameters.

The simulations clearly show that the sample size n has the most impact on the accuracy of bootstrap estimates. Once the sample size is fixed, the "goodness" of the bootstrap estimator is essentially constant independent of the number of bootstrap replications N , provided that N is above a minimum value which in this investigation turned out to be about 200.

For the estimation of a parameter θ , given an estimator $\hat{\theta} = h(X)$, the bootstrap estimate $\hat{\theta}^*$ does not seem to have an edge over the conventional estimate $\hat{\theta}$. In the linear regression simulations the bootstrap point estimates of the regression coefficients are on the average very close to their actual values and for the normal linear regression their distributions approximate the theoretical normal distributions. For the exponential distribution the simulated bootstrap estimates for the scale parameter showed an average bias which is significantly larger than the bias predicted by the theory for the maximum likelihood estimator. So for the point estimation of the parameter the maximum likelihood estimator appears to be the better estimator and the extra effort of going through the process of the bootstrap method does not bring any improvement.

For the estimation of the variance of an estimator, bootstrap in both investigated scenarios provides estimates which are very close to the theoretical results. While in the case of the exponential distribution the bootstrap estimator on the average slightly overestimates the variance, for the normal linear regression problem the bootstrap estimator is slightly below the theoretical value.

The coverage of bootstrap confidence intervals is below the nominal level for both the percentile method and the bias-corrected percentile method and not significantly different between the exponential and the linear regression cases. The latter do not show significant differences for the different distributions, viz., normal distribution or mixture of normal distributions. For small sample sizes ($n = 10$ and 20) the low coverage seems to indicate that the resulting confidence intervals are of little use. However for $n = 40$ or 50 , the coverage increases and approaches the nominal level.

An important question in practical applications is, how many bootstrap replications should be taken. The simulations show that for the estimation of the variance of the maximum likelihood estimator of the scale parameter of the exponential distribution 200 bootstrap replications are sufficient for all sample sizes. For the higher sample sizes ($n = 40, 50$) even less bootstrap replications ($N = 50, 100$) produce results of similar accuracy. For the estimation of the variance of the estimator for the coefficients in the linear regression, the simulations show that 100 bootstrap replications provide reasonable estimates. The simulations for both scenarios show that increasing the number of bootstrap replications beyond the values indicated does not significantly increase the quality of the results. For the estimation of confidence intervals the answer is not as straightforward as for the variance. If the coverage of the actual value of the parameter by the confidence interval is taken as a measure for the quality of the confidence interval estimate, the simulation results for both scenarios do not show any significant influence of the number of bootstrap replications. The investigations of the percentiles for the exponential distribution also confirm this conclusion. While Efron and Tibshirani [Ref. 7] and Efron [Ref. 12] state that for confidence intervals a minimum of 1000 bootstrap replications is required, the simulations in the special case of the exponential distribution indicate that 400 bootstrap replications would be sufficient to obtain reasonable confidence interval estimates.

APPENDIX A. PERCENTILE ESTIMATION--PERCENTILE METHOD

Table 23. 5TH PERCENTILE: Average values, standard deviations () and lengths of the 90% confidence interval (()) for bootstrap estimates of the 5th percentile of $\hat{\lambda}^*$, $\lambda = 1$.

N	n	10	20	30	40	50
200		0.7531 (0.2712) ((0.8209))	0.7789 (0.1959) ((0.6157))	0.7955 (0.1516) ((0.5052))	0.8067 (0.1332) ((0.4528))	0.8233 (0.1193) ((0.4211))
400		0.7581 (0.2769) ((0.8604))	0.7648 (0.1880) ((0.6375))	0.7937 (0.1418) ((0.4861))	0.8093 (0.1317) ((0.4644))	0.8347 (0.1255) ((0.4410))
600		0.7582 (0.2758) ((0.8291))	0.7627 (0.1792) ((0.5933))	0.7970 (0.1511) ((0.5074))	0.8299 (0.1414) ((0.4964))	0.8167 (0.1172) ((0.3813))
800		0.7620 (0.2650) ((0.8345))	0.7597 (0.1817) ((0.6060))	0.7935 (0.1388) ((0.4443))	0.8085 (0.1351) ((0.4574))	0.8361 (0.1294) ((0.4323))
1000		0.7514 (0.2708) ((0.8424))	0.7794 (0.1845) ((0.5913))	0.8013 (0.1569) ((0.5032))	0.8212 (0.1385) ((0.4659))	0.8170 (0.1348) ((0.4539))

Table 24. 95TH PERCENTILE: Average values, standard deviations () and lengths of the 90% confidence interval (()) for bootstrap estimates of the 95th percentile of $\hat{\lambda}^*$, $\lambda = 1$.

N	n	10	20	30	40	50
	200	1.979 (0.8067) ((2.424))	1.574 (0.4163) ((1.272))	1.419 (0.2754) ((0.928))	1.325 (0.2173) ((0.745))	1.290 (0.1895) ((0.656))
	400	1.996 (0.8235) ((2.501))	1.549 (0.3944) ((1.341))	1.415 (0.2626) ((0.873))	1.337 (0.2224) ((0.748))	1.315 (0.1886) ((0.670))
	600	1.999 (0.8299) ((2.625))	1.536 (0.3721) ((1.264))	1.428 (0.2831) ((0.928))	1.369 (0.2389) ((0.807))	1.292 (0.1802) ((0.644))
	800	1.996 (0.7879) ((2.568))	1.534 (0.3775) ((1.299))	1.425 (0.2681) ((0.963))	1.343 (0.2284) ((0.747))	1.311 (0.1956) ((0.678))
	1000	1.955 (0.7518) ((2.356))	1.571 (0.3740) ((1.259))	1.437 (0.2875) ((0.963))	1.368 (0.2358) ((0.780))	1.288 (0.2141) ((0.764))

APPENDIX B. PERCENTILE ESTIMATION--BIAS-CORRECTED PERCENTILE METHOD

Table 25. 5TH PERCENTILE: Average values, standard deviations () and lengths of the 90% confidence interval (()) for bootstrap estimates of the 5th percentile of $\hat{\lambda}^*$, $\lambda = 1$.

N	n	10	20	30	40	50
	200	0.7347 (0.2734) ((0.844))	0.7551 (0.2004) ((0.635))	0.7819 (0.1564) ((0.530))	0.7946 (0.1318) ((0.445))	0.8155 (0.1194) ((0.432))
	400	0.7404 (0.2736) ((0.834))	0.7444 (0.1823) ((0.615))	0.7819 (0.1520) ((0.494))	0.7979 (0.1315) ((0.477))	0.8222 (0.1214) ((0.418))
	600	0.7381 (0.2708) ((0.828))	0.7396 (0.1692) ((0.557))	0.7858 (0.1491) ((0.510))	0.8215 (0.1363) ((0.472))	0.8073 (0.1154) ((0.393))
	800	0.7524 (0.2644) ((0.817))	0.7432 (0.1751) ((0.586))	0.7755 (0.1362) ((0.422))	0.8040 (0.1306) ((0.419))	0.8297 (0.1281) ((0.435))
	1000	0.7251 (0.2586) ((0.802))	0.7680 (0.1851) ((0.591))	0.7879 (0.1581) ((0.500))	0.8117 (0.1395) ((0.461))	0.8050 (0.1299) ((0.430))

Table 26. **95TH PERCENTILE:** Average values, standard deviations () and lengths of the 90% confidence interval (()) for bootstrap estimates of the 95th percentile of $\hat{\lambda}^*$, $\lambda = 1$.

N	n	10	20	30	40	50
	200	1.888 (0.7564) ((2.293))	1.522 (0.3970) ((1.201))	1.386 (0.2683) ((0.934))	1.305 (0.2151) ((0.749))	1.273 (0.1888) ((0.653))
	400	1.900 (0.7534) ((2.321))	1.504 (0.3732) ((1.267))	1.387 (0.2569) ((0.868))	1.361 (0.2171) ((0.774))	1.297 (0.1849) ((0.658))
	600	1.908 (0.7625) ((2.476))	1.492 (0.3537) ((1.193))	1.401 (0.2733) ((0.897))	1.350 (0.2345) ((0.818))	1.276 (0.1785) ((0.638))
	800	1.911 (0.7304) ((2.350))	1.419 (0.3595) ((1.235))	1.396 (0.2574) ((0.892))	1.322 (0.2265) ((0.739))	1.259 (0.1942) ((0.665))
	1000	1.873 (0.6977) ((2.191))	1.529 (0.3568) ((1.192))	1.408 (0.2781) ((0.915))	1.345 (0.2290) ((0.764))	1.272 (0.2084) ((0.108))

APPENDIX C. VARIABILITY OF BOOTSTRAP POINT ESTIMATES--LINEAR REGRESSION

**Table 27. VARIABILITY OF THE BOOTSTRAP ESTIMATE:
Y-INTERCEPT:** Quantiles for the bootstrap estimate of the y-intercept, compared to the normal quantiles (in parentheses); theoretical value $\beta_0 = 1.5$.

Quantile n	10	20	30	40	50
0.010	0.3409 (0.3763)	0.7080 (0.7359)	0.8260 (0.8840)	0.9660 (0.9699)	1.023 (1.0277)
0.025	0.5456 (0.5532)	0.8576 (0.8562)	0.9591 (0.9810)	1.018 (1.0534)	1.062 (1.1021)
0.050	0.7194 (0.7055)	0.9476 (0.9597)	1.043 (1.0645)	1.127 (1.1252)	1.135 (1.1660)
0.100	0.8812 (0.8810)	1.063 (1.0790)	1.154 (1.1607)	1.218 (1.2080)	1.237 (1.2398)
0.250	1.173 (1.1742)	1.266 (1.2784)	1.322 (1.3214)	1.338 (1.3463)	1.368 (1.3631)
0.500	1.500 (1.5000)	1.475 (1.5000)	1.485 (1.5000)	1.498 (1.5000)	1.511 (1.5000)
0.750	1.823 (1.8258)	1.700 (1.7216)	1.665 (1.6786)	1.649 (1.6537)	1.644 (1.6269)
0.900	2.112 (2.1190)	1.903 (1.9210)	1.830 (1.8393)	1.794 (1.7920)	1.769 (1.7602)
0.950	2.275 (2.2945)	2.042 (2.0403)	1.934 (1.9355)	1.907 (1.8748)	1.845 (1.8349)
0.975	2.449 (2.4468)	2.160 (2.1438)	2.035 (2.0190)	1.982 (1.9466)	1.916 (1.8979)
0.990	2.607 (2.6237)	2.289 (2.2641)	2.119 (2.1160)	2.092 (2.0300)	1.946 (1.9723)

Table 28. VARIABILITY OF THE BOOTSTRAP ESTIMATE: SLOPE:
Quantiles for the bootstrap estimate of the slope, compared to the normal quantiles (in parentheses); theoretical value $\beta_0 = 0.8$,

Quantile n	10	20	30	40	50
0.010	0.6161 (0.6189)	0.6727 (0.6724)	0.6919 (0.6959)	0.7039 (0.7099)	0.7125 (0.7194)
0.025	0.6471 (0.6474)	0.6963 (0.6925)	0.7085 (0.7123)	0.7185 (0.7241)	0.7195 (0.7321)
0.050	0.6729 (0.6719)	0.7179 (0.7098)	0.7237 (0.7264)	0.7358 (0.7363)	0.7352 (0.7430)
0.100	0.7022 (0.7002)	0.7347 (0.7297)	0.7427 (0.7427)	0.7520 (0.7504)	0.7510 (0.7556)
0.250	0.7472 (0.7475)	0.7662 (0.7630)	0.7702 (0.7698)	0.7770 (0.7739)	0.7752 (0.7766)
0.500	0.8008 (0.8000)	0.8027 (0.8000)	0.8021 (0.8000)	0.7999 (0.8000)	0.8005 (0.8000)
0.750	0.8503 (0.8525)	0.8388 (0.8370)	0.8333 (0.8302)	0.8266 (0.8261)	0.8220 (0.8234)
0.900	0.8988 (0.8998)	0.8713 (0.8703)	0.8584 (0.8573)	0.8513 (0.8496)	0.8436 (0.8444)
0.950	0.9298 (0.9281)	0.8919 (0.8902)	0.8771 (0.8736)	0.8663 (0.8637)	0.8600 (0.8570)
0.975	0.9580 (0.9526)	0.9136 (0.9075)	0.8954 (0.8877)	0.8814 (0.8759)	0.8736 (0.8679)
0.990	0.9789 (0.9811)	0.9362 (0.9276)	0.9138 (0.9041)	0.8927 (0.8901)	0.8807 (0.8806)

APPENDIX D. FORTRAN PROGRAM FOR BOOTSTRAP

The program listed here, TEST, was written in the development of the simulations for this thesis. Its primary purpose is the validation of the simulations by repeating Efron's experiment [Ref. 4, pp. 84] as described in Chapter III. It is listed here as an example for how the bootstrap method, the percentile method and the bias-corrected percentile method were implemented in the various simulations. The program is written in FORTRAN 77 and designed to run under the IBM VM/CMS operating system.

The main program contains the generation of the original sample, the generation of the bootstrap samples and the computation of the estimator for each bootstrap sample. After sorting the bootstrap estimates the percentiles are computed with the percentile method and the bias-corrected percentile method.

The subroutine SHELLS is an implementation of the SHELL sort algorithm.

The subroutine NORMAL computes probabilities for the standard normal distribution based on the following approximation formula [Ref. 13, p. 932]

$$\Phi(z) = P(Z \leq z) =$$

$$= 1 - 0.5 \times (1 + Az + Bz^2 + Cz^3 + Dz^4 + Ez^5 + Fz^6)^{-16} \quad \text{for } z \geq 0$$

$$\text{with } A = 0.0498673470$$

$$B = 0.0211410061$$

$$C = 0.0032776263$$

$$D = 0.0000380036$$

$$E = 0.0000488906$$

$$F = 0.0000053830.$$

The subroutine INVNOR computes the quantiles for the standard normal distribution using the approximation formula [Ref. 13, p. 933]

$$z_p = \Phi^{-1}(p)$$

$$= T - \frac{A + BT + CT^2}{1 + DT + ET^2 + FT^3}$$

with

$$T = \sqrt{\ln \frac{1}{p^2}} \text{ for } p \geq 0.5$$

and

$$A = 2.515517$$

$$B = 0.802853$$

$$C = 0.010328$$

$$D = 1.432788$$

$$E = 0.189269$$

$$F = 0.001308.$$

The subroutines NORMAL and INVNOR are both used for the bias-corrected percentile method.

PROGRAM TEST

```

*****
* Program to verify implementation of methods by repeating Efron's
* simulation for comparison.
*
* Simulate the bootstrap from an Exponential distribution with
* parameter LAMBDA with standardized original samples.
* Compute 5th, 10th, 90th and 95th percentiles using the
*   - Percentile Method
*   - Bias-Corrected Percentile Method
*
* Variables and parameters:
*   N   Original sample size = 15
*   NN  Number of bootstrap replications = 1000
*   M   Number of repetitions = 10
*   ORIG the original sample
*   RAND the Uniform (0,1) random numbers for the bootstrap
*   DRAW the integer random numbers for the bootstrap
*   LHAT the vector of MLEs of the bootstraps
*****

* Declare variables and I/O devices

INTEGER N, NN, M, IX1, IX2, ISORT, MUL1, MUL2, DRAW, LOOK
PARAMETER (N=15, NN=1000, M=10)
REAL LAMBDA, AV5, AV10, AV90, AV95, CDFLHA, ZPRIME, Z5, Z10, Z90,
CZ95, AUX5, AUX10, AUX90, AUX95, BAV5, BAV10, BAV90, BAV95,
CAAA5, AAA10, AAA90, AAA95
PARAMETER (LAMBDA=1.0)
REAL ORIG(N), RAND(N), LHAT(NN), P5(M), P10(M), P90(M), P95(M),
CBCP5(M), BCP10(M), BCP90(M), BCP95(M)
DATA IX1/31397/, IX2/75931/, MUL1/1/, MUL2/2/, ISORT/0/, AV5/0/,
CAV10/0/, AV90/0/, AV95/0/, BAV5/0/, BAV10/0/, BAV90/0/, BAV95/0/

CALL EXCMS('FILEDEF 10 DISK OUTEST LISTING A')

* Output header and compute constants

WRITE(10,90) N, NN
90  FORMAT('1/'0',10X,'BOOTSTRAP SIMULATION/'0',10X,
C'Nonparametric confidence intervals for the expectation,'/11X,
C'negative exponential distribution;' /11X,
C'standardized samples of size n = ',I4/11X,
C'number of bootstrap replications N = ',I6/'0',4X,
C'Trial          Percentile Meth.          Bias-corr. Percentile '
C,'Meth. '/5X,
C'      5%      10%      90%      95%      5%      10%      90% '
C,' 95%/'0')

CALL INVNOR(0.95,Z5)
CALL INVNOR(0.9,Z10)
CALL INVNOR(0.1,Z90)
CALL INVNOR(0.05,Z95)

* For M repetitions

```

```

      DO 30, K = 1, M

** Create the original sample

      CALL LEXPN(IX1,ORIG,N,MUL1,ISORT)

** Standardize the original sample

      A = 0
      B = 0
      DO 10, JJ = 1, N
        A = A + ORIG(JJ)
        B = B + ORIG(JJ)**2
10      CONTINUE
      SD = SQRT((B - A*A/N)/(N-1))
      DO 11, JJ = 1, N
        ORIG(JJ) = (ORIG(JJ) - A/N)/SD
11      CONTINUE

** Do NN bootstrap replications

      DO 20, I = 1, NN

        CALL LRND(IX2,RAND,N,MUL2,ISORT)
        LHAT(I) = 0
        DO 21, J = 1, N
          DRAW = INT(N*RAND(J)) + 1
          LHAT(I) = LHAT(I) + ORIG(DRAW)/N
21        CONTINUE
20      CONTINUE

** Sort the bootstrap estimates

      CALL SHELLS(LHAT,NN)

** Compute percentiles using the percentile method

      P5(K) = LHAT(50)
      P10(K) = LHAT(100)
      P90(K) = LHAT(900)
      P95(K) = LHAT(950)

** Compute the percentiles using the bias-corrcted percentile method

      LOOK = NN/2
111    IF(LHAT(LOOK).GT.0.AND.LHAT(LOOK+1).GT.0) THEN
        LOOK = LOOK - 1
        GO TO 111
      ELSE IF(LHAT(LOOK).LT.0.AND.LHAT(LOOK+1).LT.0) THEN
        LOOK = LOOK + 1
        GO TO 111
      END IF

      CDFLHA = REAL(LOOK)/NN
      CALL INVNOR(CDFLHA,ZPRIME)

```

```

AUX5 = 2*ZPRIME - Z5
AUX10 = 2*ZPRIME - Z10
AUX90 = 2*ZPRIME - Z90
AUX95 = 2*ZPRIME - Z95
CALL NORMAL(AUX5,AAA5)
CALL NORMAL(AUX10,AAA10)
CALL NORMAL(AUX90,AAA90)
CALL NORMAL(AUX95,AAA95)
BCP5(K) = LHAT( INT( AAA5*NN) )
BCP10(K) = LHAT( INT( AAA10*NN) )
BCP90(K) = LHAT( INT( AAA90*NN) )
BCP95(K) = LHAT( INT( AAA95*NN) )

```

** Output trial results

```

      WRITE(10,91) K,P5(K), P10(K), P90(K), P95(K), BCP5(K),
C      BCP10(K), BCP90(K), BCP95(K)
91      FORMAT('0',I8,4(2X,F6.3),3X,4(2X,F6.3))
30      CONTINUE

```

* Compute and output averages

```

      DO 55, K = 1, M
        AV5 = AV5 + P5(K)/M
        AV10 = AV10 + P10(K)/M
        AV90 = AV90 + P90(K)/M
        AV95 = AV95 + P95(K)/M
        BAV5 = BAV5 + BCP5(K)/M
        BAV10 = BAV10 + BCP10(K)/M
        BAV90 = BAV90 + BCP90(K)/M
        BAV95 = BAV95 + BCP95(K)/M
55      CONTINUE
      WRITE(10,92) AV5, AV10, AV90, AV95, BAV5, BAV10, BAV90, BAV95
92      FORMAT('0',1X,'Average',4(2X,F6.3),3X,4(2X,F6.3))
      STOP
      END

```

SUBROUTINE SHELLS(UNSORT,NUM)

* Subroutine SHELLS to sort data in ascending order (Shell-sort) *

```

      INTEGER NUM, GAP, COUNT
      REAL UNSORT(NUM)

      GAP = NUM
10      GAP = INT(GAP/2.0)
20      COUNT = 0
      DO 40, I = 1, NUM - GAP
        IF (UNSORT(I).LE.UNSORT(I+GAP)) GO TO 40
        A = UNSORT(I)
        UNSORT(I) = UNSORT(I+GAP)

```



```

        UNSORT(I+GAP) = A
        COUNT = COUNT + 1
40    CONTINUE
        IF (COUNT.GT.0) GO TO 20
        IF (GAP.GT.1) GO TO 10
        RETURN
    END

```

SUBROUTINE NORMAL(INPUT,RESULT)

```

*****
* Subroutine to compute probabilities for the standard normal      *
* distribution.                                                    *
*****

```

* Declare variables

```

    LOGICAL NEG
    DOUBLE PRECISION AA, BB, CC, DD, EE, FF, Z, X
    REAL RESULT, INPUT
    DATA AA/0.049867347D0/, BB/0.0211410061D0/, CC/0.0032776263D0/,
    CDD/0.0000380036D0/, EE/0.0000488906D0/, FF/0.000005383D0/

```

* Prepare input

```

    NEG = .FALSE.
    IF(INPUT.LT.0) NEG = .TRUE.
    Z = DBLE(INPUT)
    IF(NEG.EQV..TRUE.) Z = -Z

```

* Apply formula

```

    X = 1D0 + AA*Z + BB*Z**2D0 + CC*Z**3D0 + DD*Z**4D0
    X = X + EE*Z**5D0 + FF*Z**6D0
    X = X**(-16D0)
    X = 1D0 - 0.5D0*X

```

* Prepare output

```

    RESULT = REAL(X)
    IF(NEG.EQV..TRUE.) RESULT = 1 - RESULT
    RETURN
    END

```

SUBROUTINE INVNOR(INPUT,RESULT)

```

*****
* Subroutine to compute quantiles of the standard normal distribution. *
*****

```

* Declare variables

```
DOUBLE PRECISION A, B, C, D, E, F, P, T, YY
REAL INPUT, RESULT
LOGICAL LESS
DATA A/2.515517D0/, B/0.802853D0/, C/0.010328D0/, D/1.432788D0/,
CE/0.189269D0/, F/0.001308D0/
```

* Prepare input

```
LESS = .FALSE.
P = DBLE(INPUT)
IF(INPUT.LT.0.5) LESS = .TRUE.
IF(LESS.EQV..FALSE.) P = 1D0 - P
```

* Apply formula

```
T = DSQRT(DLOG(1D0/P**2D0))
YY = A + B*T + C*T**2D0
YY = YY/(1D0 + D*T + E*T**2D0 + F*T**3D0)
YY = T - YY
```

* Prepare result

```
IF(LESS.EQV..TRUE.) YY = -YY
RESULT = REAL(YY)
RETURN
END
```

APPENDIX E. SIMTBED DRIVER FOR BOOTSTRAP

The program listed here is an example for the drivers used in the simulations under SIMTBED [Ref. 9]. SIMTBED is written in FORTRAN and operates under the IBM Professional FORTRAN, which is a prerequisite for any application. The drivers, i. e. the user input has to be written in this programming language.

The main part of the driver is the general call to SIMTBED with all the parameters like sample size, number of replications, destination and form of the output etc.. This part basically follows the instructions in the manual.

The subroutines constitute the part of the driver, which is specific for each problem. Here the user has to set up the specific simulation by choosing the distribution of the random variates and by programming how the statistical estimates are to be computed. The subroutines LIRE1 and LIRE2 in the driver are written to compute the bootstrap estimate of the y-intercept respectively the slope in normal linear regression. In all simulations involving normal linear regression the same basic setup was used.

```

*****
*      SIMTBED driver for normal linear regression      *
*      Main program                                   *
*****

```

* Declare variables

```

      REAL*4  Y(3500),YMIN, YMAX
      CHARACTER*120 T1,T2
      REAL*8  IX1,IX2,IX3,IX4,IX5,IX
      INTEGER  N,M,NE(8),L,D,RG,SEI,SVS,NEST,NCOLRNDX(3),IFILE,IBWPRT,
C MSE, NPRT, IPR, IBIV, IRSTR
      REAL VMSE(8,5),VMX1(8,4),VMX2(8,4),VMX3(8,4),VMX4(8,4),VMX5(8,4)
      EXTERNAL LIRE1, LIRE2

```

* Input of SIMTBED parameters

```

      DATA  N/3000/
      DATA  M/ 1/
      DATA  NE/ 10,20,30,40,50,50,50,50/
      DATA  L/5/
      DATA  D/ 0/
      DATA  RG/ 0/
      DATA  SEI/ 0/
      DATA  SVS/ 0/
      DATA  YMIN/ 0. /
      DATA  YMAX/ 0. /
      DATA  IX/ 88771.DO/
      DATA  IFILE /0/
      DATA  NPRT  /1/
      DATA  MSE    /0/
      DATA  VMSE  /40*0/
      DATA  IPR    /0/
      DATA  IBIV   /0/
      DATA  IRSTR  /1/
      DATA  ICOLOR/0/, IBWPRT/1/, NCPRT/1/, NCOLRNDX/1,2,7/

```

* Set output parameters

```

      DATA  T1/'Y-INTERSECT NORMAL LINEAR REGRESSION
C (100 Bootstrap replications)'/
      DATA  T2/'SLOPE  NORMAL LINEAR REGRESSION
C (100 Bootstrap replications)'/

      OPEN(06,FILE='LIRE100. OUT',ERR=999,IOSTAT=IER)
      OPEN(05,FILE='CON',ERR=999,IOSTAT=IER)
      OPEN(02,FILE='LIRE100. RST',ERR=999,IOSTAT=IER,FORM='UNFORMATTED',
C ACCESS='SEQUENTIAL')
      OPEN(01,FILE='LIRE100. DAT',ERR=999,IOSTAT=IER,FORM='FORMATTED',
C ACCESS='SEQUENTIAL')

```

* Generator parameters

```

      NEST=2
      NSR=10

```

```

IX1=IX
IX2=IX
IX3=IX
IX4=IX
IX5=IX

```

```
* Make the call to SIMTBED
```

```

CALL SMTBED(IX1,IX2,IX3,IX4,IX5,Y,N,M,NE,L,D,NSR,RG,SEI,SVS,
C YMIN,YMAX,NEST,LIRE1,T1,LIRE2,T2,LIRE1,T1,LIRE1,T1,LIRE1,T1,
C IFILE,NPRT,MES,VMSE,IPR,VMX1,VMX2,VMX3,VMX4,VMX5,IBIV,IRSTR,
C ICOLOR,IBWPRT,NCPRT,NCOLRNDX)

```

```
STOP
```

```

999 CONTINUE
WRITE(6,*) 'ERROR OPENING FILE 1, 2 OR 6'
END

```

```
*****
```

```
SUBROUTINE LIRE1(ISEED,N,EVAL)
```

```

*****
* Subroutine to evaluate the first estimator, the Y-intercept *
*****

```

```
* Declare variables
```

```

INTEGER N, BREP, DRAW, NN
REAL BETA0, BETA1, B0, B1, BHAT0, BHAT1, VAR, XBAR, YBAR,
C NUM1, DENOM, EVAL
REAL*8 ISEED
PARAMETER (VAR=0.5, BETA0=1.5, BETA1=0.8, NN=50)
REAL X(NN), Y(NN), EPSI(NN), RAND(NN), NORRAN(NN)

```

```
BREP = 100
```

```
* Compute the x-values and related results
```

```

DO 99, I = 1, N
  X(I) = REAL(I)*10./REAL(N)
99 CONTINUE

XBAR = 0
DO 1, I = 1, N
  XBAR = XBAR + X(I)/N
1 CONTINUE

DENOM = 0
DO 2, I = 1, N
  DENOM = DENOM + (X(I) - XBAR)**2
2 CONTINUE

```

* Create original pairs of observations and compute parameters

```
      CALL LNORPC(ISEED,NORRAN,N)

      YBAR = 0
      DO 10, I = 1, N
        Y(I) = BETA0 + BETA1*X(I) + SQRT(VAR)*NORRAN(I)
        YBAR = YBAR + Y(I)/N
10     CONTINUE

      NUM1 = 0
      DO 11, I = 1, N
        NUM1 = NUM1 + (X(I) - XBAR)*(Y(I) - YBAR)
11     CONTINUE

      B1 = NUM1/DENOM
      B0 = YBAR - B1*XBAR
```

* Compute the epsilons

```
      DO 12, I = 1, N
        EPSI(I) = Y(I) - B0 - B1*X(I)
12     CONTINUE
```

* Do the bootstraps

```
      BOBAR = 0
      B1BAR = 0
      DO 20, J = 1, BREP
        CALL LRNDPC( ISEED,RAND,N)
        YBAR = 0
        NUM1 = 0
        DO 21, K = 1, N
          DRAW = INT(N*RAND(K)) + 1
          Y(K) = B0 + B1*X(K) + EPSI(DRAW)
          YBAR = YBAR + Y(K)/REAL(N)
21       CONTINUE

        DO 22, K = 1, N
          NUM1 = NUM1 + (Y(K) - YBAR)*(X(K) - XBAR)
22       CONTINUE

        B1HAT = NUM1/DENOM
        BOBAR = BOBAR + (YBAR - B1HAT*XBAR)/REAL(BREP)
20      CONTINUE

      EVAL = BOBAR

      RETURN
      END
```

SUBROUTINE LIRE2(ISEED,N,EVAL)

```
*****
* Subroutine to evaluate the second estimator, the Slope *
*****
```

* Declare variables

```
INTEGER N, BREP, DRAW, NN
REAL BETA0, BETA1, B0, B1, BHAT0, BHAT1, VAR, XBAR, YBAR,
C NUM1, DENOM, EVAL
REAL*8 ISEED
PARAMETER (VAR=0.5, BETA0=1.5, BETA1=0.8, NN=50)
REAL X(NN), Y(NN), EPSI(NN), RAND(NN), NORRAN(NN)
```

BREP = 100

* Compute the x-values and related results

```
DO 99, I = 1, N
  X(I) = REAL(I)*10./REAL(N)
99  CONTINUE

XBAR = 0
DO 1, I = 1, N
  XBAR = XBAR + X(I)/N
1  CONTINUE

DENOM = 0
DO 2, I = 1, N
  DENOM = DENOM + (X(I) - XBAR)**2
2  CONTINUE
```

* Create original pairs of observations and compute parameters

```
CALL LNORPC(ISEED,NORRAN,N)

YBAR = 0
DO 10, I = 1, N
  Y(I) = BETA0 + BETA1*X(I) + SQRT(VAR)*NORRAN(I)
  YBAR = YBAR + Y(I)/N
10  CONTINUE

NUM1 = 0
DO 11, I = 1, N
  NUM1 = NUM1 + (X(I) - XBAR)*(Y(I) - YBAR)
11  CONTINUE

B1 = NUM1/DENOM
B0 = YBAR - B1*XBAR
```

* Compute the epsilons

```
DO 12, I = 1, N
```

```

      EPSI(I) = Y(I) - B0 - B1*X(I)
12    CONTINUE

* Do the bootstraps

      BOBAR = 0
      B1BAR = 0
      DO 20, J = 1, BREP
        CALL LRNDPC(ISEED,RAND,N)
        YBAR = 0
        NUM1 = 0
        DO 21, K = 1, N
          DRAW = INT(N*RAND(K)) + 1
          Y(K) = B0 + B1*X(K) + EPSI(DRAW)
          YBAR = YBAR + Y(K)/REAL(N)
21      CONTINUE

        DO 22, K = 1, N
          NUM1 = NUM1 + (Y(K) - YBAR)*(X(K) - XBAR)
22      CONTINUE

        B1HAT = NUM1/DENOM
        B1BAR = B1BAR + B1HAT/REAL(BREP)
20    CONTINUE

      EVAL = B1BAR

      RETURN
      END

```

LIST OF REFERENCES

1. Efron, Bradley, *Bootstrap Methods: Another Look at the Jackknife*, The Annals of Statistics, 1979, Vol. 7, No. 1, pp. 1 - 26.
2. Beran, Rudolf, *Estimated Sampling Distributions: The Bootstrap and Competitors*, The Annals of Statistics, 1982, Vol. 10, No. 1, pp. 212 - 225.
3. Bickel, Peter J. and Freedman, David A., *Some Asymptotic Theory for the Bootstrap*, The Annals of Statistics, 1981, Vol. 9, No. 6, pp. 1196 - 1217.
4. Efron, Bradley, *The Jackknife, the Bootstrap and Other Resampling Plans*, Society for Industrial and Applied Mathematics, 1982.
5. Davison, A. C., Hinkley, D. V. and Schechtman, E., *Efficient Bootstrap Simulation*, Biometrika, 1986, Vol. 73, No. 3, pp. 555 - 566.
6. Schenker, Nathaniel, *Qualms About Bootstrap Confidence Intervals*, Journal of the American Statistical Association, June 1985, Vol. 80, No. 390, pp. 360 - 361.
7. Efron, B. and Tibshirani, R., *Bootstrap Methods for Standard Errors, Confidence Intervals, and Other Measures of Statistical Accuracy*, Statistical Science, 1986, Vol. 1, No. 1, pp. 54 - 75.
8. Cortes-Colon, William, *An Analysis of the Bootstrap Method for Estimating the Mean Squared Error of Statistical Estimators*, M. S. Thesis, Naval Postgraduate School, Monterey, California, September 1986.
9. Lewis, P. A. W., Orav, E. J., and Uribe, L., *Advanced Simulation and Statistics Package*, Wardworth and Brooks, 1986, with most recent changes by Youmans, R. in his concurrent thesis, Naval Postgraduate School, Monterey, September 1988.

10. *GRAFSTAT: An APL System for Interactive Scientific-Engineering Plotting, Data Analysis, Applied Statistics, and Graphics Output Development*, IBM Research.
11. Naval Postgraduate School Report NPS55-81-005, *The New Naval Postgraduate School Random Number Package LLRANDOMII*, by P. A. W. Lewis and L. Uribe, February 1981.
12. Efron, Bradley, *Better Bootstrap Confidence Intervals*, Journal of the American Statistical Association, March 1987, Vol. 82, No. 397, pp. 171 - 185.
13. Abramowitz, M. and Stegun, I. A., *Handbook of Mathematical Functions*, National Bureau of Standards, Applied Mathematics Series 55.

INITIAL DISTRIBUTION LIST

	No. Copies
1. Defense Technical Information Center Cameron Station Alexandria, VA 22304-6145	2
2. Library, Code 0142 Naval Postgraduate School Monterey, CA 93943-5002	2
3. Prof. Toke Jayachandran Naval Postgraduate School (Code 53Jy) Department of Mathematics Monterey, California 93943-5000	2
4. Prof. Robert R. Read Naval Postgraduate School (Code 55Re) Department of Operations Research Monterey, California 93943-5000	2
5. Bundesminister der Verteidigung Fü H IV 1 Postfach 1378 5300 Bonn 1 West Germany	1
6. Amt für Studien und Übungen der Bundeswehr Friedrich-Ebert-Str. 72 5660 Bergisch Gladbach 1 West Germany	1
7. Amt für Studien und Übungen der Bundeswehr Mil. Bereich OR Einstein-Str. 20 8012 Ottobrunn West Germany	1
8. Hptm. Stefan Bernhardt 4./Instandsetzungsbataillon 110 4420 Coesfeld West Germany	1

Thesis
B452851 Bernhardt
c.1 Small sample proper-
ties of bootstrap.

Thesis
B452851 Bernhardt
c.1 Small sample proper-
ties of bootstrap.

thesB452851

Small sample properties of bootstrap.



3 2768 000 84759 4

DUDLEY KNOX LIBRARY